Snowflake BDT

Jan Kozák

22 October 2025

Another technology. Why?

- > Technologies just evolve and change ... nowadays maybe even faster
- 2000s RDBMS, vertical scaling
 - Relational data, DWH
- 2010s clusters, horizontal scaling
 - Data lake concept, unstructured and semistructured data
 - BDT begins on-premise Hadoop cluster, MapReduce / Spark
 - Open Source ... but still may lead to vendor lock-in (Cloudera)
- 2020s widespread cloud adoption
 - Lakehouse, unified data platforms ...
 - Data for AI
- 2030s ???

What is Snowflake





THE SNOWFLAKE PLATFORM

"A single, fully managed platform that powers the AI Data Cloud. Snowflake securely connects businesses globally across any type or scale of data to productize AI, applications and more in the enterprise."

https://www.snowflake.com/en/data-cloud/platform/



Brief History

- Founded in 2012 by former Oracle engineers Benoit Dageville, Thierry Cruanes
- Relational DB in cloud initially focused on DWH
 - Resource elasticity (data, compute)
 - Fault isolation
 - Performance isolation
- > Currently
 - 10,000+ customers
 - 5 billion queries processed daily
 - 2,500 Marketplace listings

Core concepts

- Cloud agnostic
 - AWS, Azure, GCP
 - Interoperability
 - Between SNFL accounts
 - Between SNFL account and various cloud providers

Cloud oriented

- Data ingestion primary sources are data stored in cloud services (Azure Blob Storage, AWS S3, Google Cloud Storage)
- Cloud only, infrastructure provisioned by Snowflake
 - You cannot deploy Snowflake on your local device
 - Neither BYOC

Account

PROFINIT >

- Organization company or similar
 - Accounts linked with organization
- Account
 - Snowflake account ~ DB instance/machine
 - Cloud provider, region cannot change later
 - Edition (Standard, Enterprise, Business Critical, Virtual Private Snowflake)
- > Want to try it out?
 - https://signup.snowflake.com/
 - 30 days for free, 400 USD credit
 - Select at least Enterprise edition to unlock most features
 - So far no restrictions on number of trial accounts you may create more than one at time or register a new one when the trial period expires

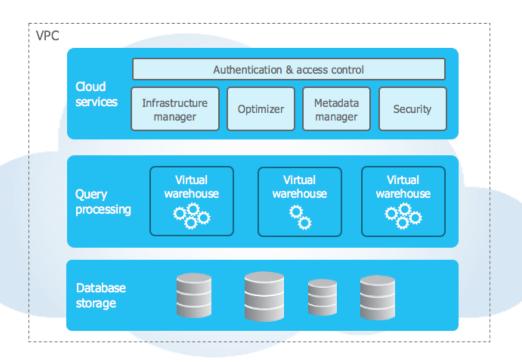
Snowflake - Architecture

Architecture

- Self-managed Data Platform
 - Snowflake handles all infrastructure (HW, SW, maintenance)
- Data storage and compute resources separation
 - Compute Layer (Virtual Warehouse, Serverless)
 - Data Storage
 - Cloud Services
- > Hybrid of shared-disk and shared-nothing architectures
 - Shared-disk data stored at a central data repository accessible from all compute nodes
 - Shared-nothing MPP each node in the cluster stores a portion of data set locally during computation

Architecture





Ref: https://docs.snowflake.com/en/user-guide/intro-key-concepts

Compute layer

- Virtual Warehouse (WH)
 - Main computing resource
 - Active WH is necessary for most operations
- May create multiple WH per account
 - Separate monitoring, access controls, sizing, tagging ...
- Sizing
 - Currently 10 sizes from X-Small (XS) to 6X-Large
 - + Snowpark Optimized counterparts for M-6XL sizes ML/AI tasks
 - Exact technical specifications are not public
- Scaling both up/down and in/out (even online)
 - Scaling-out only in Enterprise and higher editions, 1-10 clusters

Compute layer II

- Serverless cannot create directly, only for tasks/pipes (see later)
- > What if I need GPUs?
 - Last year's addition to Snowflake
 - Snowpark Container Services
 - Analogy to Docker / Kubernetes / OCP
 - Image run on a compute pool instances

Cost attribution

- I know who ran what workload and how many credits it consumed
- Great leap from RDBMS

So what is in Snowflake?





What is Snowflake made of?

- Database objects
 - Everything is a database object
 - Database, Schema, Table, Views, ...
 - Users, Roles, Tasks, Integrations, Notebooks, ...
- > SQL command for everything
 - Connection to git? Set up API integration

```
CREATE OR REPLACE API INTEGRATION my_git_api_integration

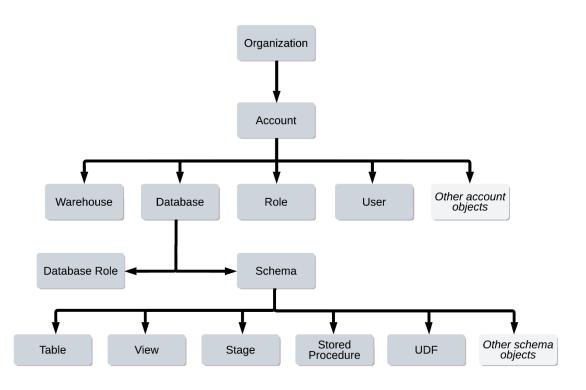
API_PROVIDER = git_https_api

API_ALLOWED_PREFIXES = ('https://github.com/my-account')

ALLOWED_AUTHENTICATION_SECRETS = (my_git_secret)

ENABLED = TRUE;
```

Objects Hierarchy



Ref: https://docs.snowflake.com/en/user-guide/security-access-control-overview

Database Objects

- Structure fairly similar to common relational DBs
 - Database schema tables/views/functions...
-) Differences
 - Role (user/account-level) vs Database Role
 - RBAC Snowflake cares about roles, not about users
 - Stage landing point for data (un)loading
 - Internal primarily for ad-hoc data uploaded from local machines
 - External connection to cloud data storage accounts
 - AWS S3, Azure Blob Storage, Google Cloud Storage
 - Other special types
 - Integrations, Masking Policies, Data Quality Measures, ...

Storing the data

- Many table types serving different purposes/use-cases
 - Internal Table
 - Transient, Temporary
 - External Table
 - Directory Table
 - Dynamic Table
 - Event Table
 - Hybrid Table
- Main difference to standard relational DBs:
 - PKs, FKs unenforced
 - No indexes, no (explicit) partitioning

Data Storage

- Using data storage services of underlying cloud providers
 - VPS storage account dedicated only for the account
 - Others all accounts in one region share one centralized storage account
- Data Format (for internal tables)
 - Column-oriented
 - Proprietary and private
- Micro-partitioning
 - Automatically for all Snowflake internal tables
 - Compressed data file (disk segment), typically 50-500 MB (pre-compression)
 - Metadata about value ranges, selectivity etc. -> efficient filtering and pruning
 - Snowflake claims it is efficient for tables up to 5 TB large

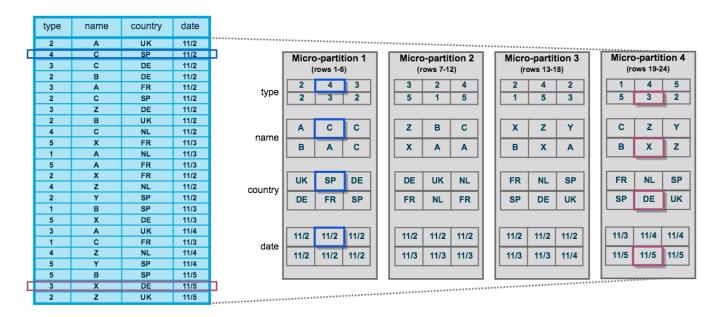
Micro-partitions



Table: ±1

Logical Structure

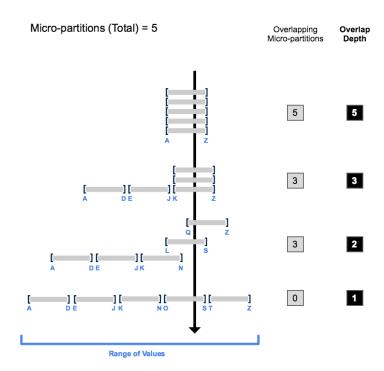
Physical Structure



Ref: https://docs.snowflake.com/en/user-guide/tables-clustering-micropartitions

Clustering depth

- Measure for efficient query pruning
- May be different for each column
- Lower depth ~ faster queries
- Affected by load patterns
 - Random vs bulk load



Ref: docs.snowflake.com/en/user-guide/tables-clustering-micropartitions

Clustering, Other data formats

- Micro-partitioning is used in internal tables
- When micro-partitioning is not performing well, you may declare clustering keys
- Support for unstructured and semi-structured data
 - Semi-structured
 - JSON, XML, Avro, Parquet, ORC directly supported by native functions
 - VARIANT column
 - Unstructured
 - BLOB/CLOB very limited size in internal tables (16 MB)

Table Features

- Time travel
 - pros: error rollback, quasi-historization
 - cons: cost

```
SELECT * FROM my_table AT(TIMESTAMP => 'Wed, 26 Jun 2024 09:20:00 -0700'::timestamp_tz);

SELECT * FROM my_table AT(OFFSET => -60*5);
```

Fail-safe period

- 7-days, not for transient and temporary tables
- "Best-effort" if some critical damage happens to your data
 - Exterme operational failures
 - Security breach

Tables, views

- All table types can be queried together
 - We can join internal, external, directory and hybrid table together in one query
 - We can declare views over any table type

> Views

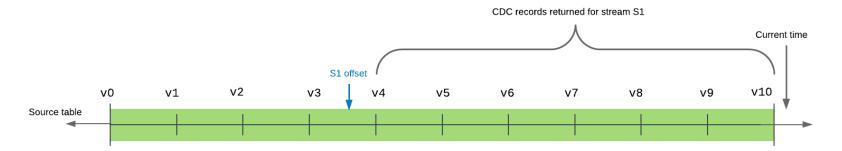
- Secure Views
 - Suppressing some internal optimizations so the consumer cannot indirectly see or deduce data they aren't privileged to (e.g. through query plan or aggregation queries)
 - Only view owners can see its definition
- Materialized Views
 - For large or complex queries, results refreshed automatically
 - Only in Enterprise and higher editions; can be also Secure Materialized Views

Data Loading and Transformations

- Relevant objects
 - Orchestration and scheduling
 - Tasks
 - Streams
 - Pipes
 - Incoming data storage
 - Stage (internal or external)
 - Table types
 - Dynamic tables
 - External tables
 - Directory tables
 - (and of course "standard" internal tables)

Batch loading

- Tasks + streams
- Stream Change Data Capture (CDC)



Consumed only once -> multiple consumers need multiple streams

Batch loading

- Task
 - Scheduled (e.g. every 10 minutes)
 - Triggered by stream or by another task
 - Serverless or running on user managed WH

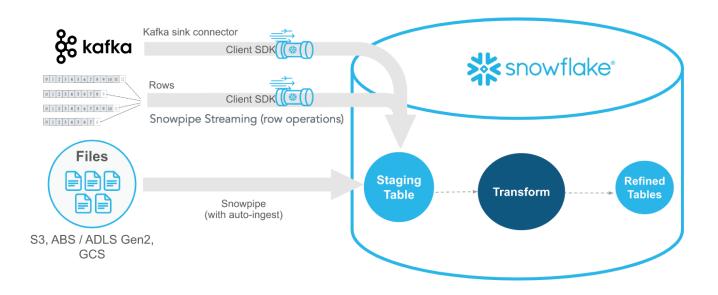
```
CREATE TASK my_triggered_task

TARGET_COMPLETION_INTERVAL='15 MINUTES'
WHEN SYSTEM$STREAM_HAS_DATA('my_order_stream')
AS

INSERT INTO CUSTOMEr_activity
SELECT CUSTOMEr_id, order_total, order_date, 'order'
FROM my_order_stream;
```

Continuous loading

- > Snowpipe
 - Monitor external stage and trigger pipeline when new data arrives
 - Map external storage to ext. stage (Amazon S3, Azure Blob Storage)
 - Consume events about new data from PubSub services
 - Azure EventGrid, Amazon SNS/SQS
- Snowpipe Streaming
 - connecting directly to external stream (Kafka)



Snowpipe Streaming | Snowflake Documentation

Dynamic tables

- Declarative transformation framework (cf. Delta Live Tables in DBX)
 - Similar results could be achieved with other previously mentioned methods, but with much more code
- Configurable "data freshness" / target lag

```
CREATE OR REPLACE DYNAMIC TABLE names

TARGET_LAG = '1 minute'

WAREHOUSE = mywh

AS

SELECT

var:id::int id,
var:fname::string first_name,
var:lname::string last_name

FROM raw;
```

Hybrid tables

- Getting back to RDBMS times
 - ACID, referential integrity
 - PKs, FKs enforced, good old indexes
- For real-time transactional use-cases
 - Random read/writes
 - Unistore concept

Functions, Procedures, Scripting

- Snowflake provides rich set of predefined functions
 - Window functions
 - Analytic functions (time-series, MATCH)
 - Approximate functions (similarity, percentiles ...)
- Viser defined functions and procedures
 - Supports multiple languages
 - SQL, Python, Java, JavaScript
- External functions (executed outside Snowflake)
- Scripting REST API, Python API
- Snowpark for data processing, similar to Apache Spark (based on DataFrame paradigm)

Data sharing

- Zero-copy cloning
- > Data replication
 - Same organization, to accounts in different regions

> Data sharing

- One of the strongest Snowflake features
- Share data with accounts from other organizations (in the same region)
- If replicated, may share data even across different regions and cloud providers
- Consumer read-only, but no additional costs (data stored only at provider)
- Multiple methods
 - Marketplace listings, Secure Data Sharing, Data Exchange, Clean Rooms

Features Overview

- Masking policies (row-level, column-level)
- > Snowflake/Snowpark ML
 - Set of services for ML and MLOps, mimicking Mlflow
 - Model Registry, Feature Store, Datasets
- > Applications
 - Streamlit, Snowflake Native App
- Snowpark Container Services
 - For Native Apps or Snowpark ML model serving

Snowflake ML/AI

- Rapidly evolving in last years/months
- 2024 introducing interactive notebooks
- > LLM Snowflake Cortex
 - Multiple models supported (mistral, llama, Snowflake Arctic ...)
 - RAG native support for vector embeddings storage (VECTOR datatype)
 - Cortex Search out-of-box RAG over stored documents
 - Cortex Analyst answering nature language questions about data
 - Document AI extraction and classification of documents

Resources - Billing

- So what does Snowflake charge you for?
 - Data storage (20-40 USD/TB per month)
 - Compute costs in term of *credits*
 - 1 Snowflake credit ~ 2-10 USD (varies by provider, region and edition)
 - Virtual warehouses 1 credit = 1 hour of XS warehouse uptime
 - Serverless actual compute time
 - Cloud services usually negligible
 - Data transfer costs
 - Egress (outbound traffic) from Snowflake account both unloading to data storage (e.g. AWS S3) or Snowflake accounts in another regions
 - Ingress generally not charged (by Snowflake, cloud providers will charge some fee...)

Resources - Control

- > Almost unlimited resources
- > If you pay for it



Resources - Monitoring

- > How to control our bills?
- > Resource Monitor
 - Formally just another DB object type
 - Controls credit consumption of virtual warehouses
 - Actions Notify, Suspend, Suspend immediately
 - But it doesn't monitor serverless features or cloud services usage
- Internal metadata tables
- Cost attribution
 - Using ACCOUNT_USAGE tables and tags
 - Per organizational unit, user or query

Polaris, Horizon



> Snowflake Horizon Catalog

- Counterpart to Unity Catalog, Snowflake governance solution
- Governing data, apps and models across whole organization in single place
- Focused on Snowflake, but may integrate (on some level) data from external sources, mainly if they use Apache Iceberg format

> Polaris Catalog

- Open source (one of first Snowflake open-sourced projects)
- Using Apache Iceberg REST API for integrations
- Deployment: not necessarily on Snowflake
- Ambition: cross-engine universal catalog, industry standard
 - Integrations with Databricks, Fabric?

Almost finishing ... some of Snowflake Tops

- Modern approach to old relational databases
 - Easy function/procedures definition in Python, JavaScript and Java
 - GitOps, dbt
 - Non-standard SQL features
 - CREATE OR ALTER
 - GROUP BY ALL, UNION BY NAME
 - Trailing comma
 - Calculated column reference by name in the same SELECT
 - And many others ...
 - Constant improvement
 - Materialized views pricing
 - Non-critical tasks scheduling

Snowflake Ecosystem





PROFINIT >

Questions



Thank you

Profinit EU, s.r.o. Tychonova 2, 160 00 Praha 6

Tel.: + 420 224 316 016, web: www.profinit.eu

in LINKEDIN linkedin.com/company/profinit





